

# Census Income Project

The dataset provided to us contains 32560 rows, and 14 different independent features. We aim to predict if a person earns more than 50k\$ per year or not. Since the data predicts 2 values (>50K or <=50K), this clearly is a classification problem, and we will train the classification models to predict the desired outputs.

## Importing Libraries

```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

## Loading Dataset

```
In [1]: import pandas as pd
df=pd.read_csv(r"https://raw.githubusercontent.com/dsrscentist/dataset1/master/census_income.csv")
df.head()
```

50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	0	13	United-States	<=50K
1	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K
2	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K
3	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K
4	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	<=50K
df															
	Age	Workclass	Fnlwgt	Education	Education_num	Marital_status	Occupation	Relationship	Race	Sex	Capital_gain	Capital_loss	Hours_per_week	Native_country	Income

```
In [18]: df
```

2	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K
3	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K
4	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	<=50K
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
32555	27	Private	257302	Assoc-acdm	12	Married-civ-spouse	Tech-support	Wife	White	Female	0	0	38	United-States	<=50K
32556	40	Private	154374	HS-grad	9	Married-civ-spouse	Machine-op-inspct	Husband	White	Male	0	0	40	United-States	>50K
32557	58	Private	151910	HS-grad	9	Widowed	Adm-clerical	Unmarried	White	Female	0	0	40	United-States	<=50K
32558	22	Private	201490	HS-grad	9	Never-married	Adm-clerical	Own-child	White	Male	0	0	20	United-States	<=50K
32559	52	Self-emp-inc	287927	HS-grad	9	Married-civ-spouse	Exec-managerial	Wife	White	Female	15024	0	40	United-States	>50K

32560 rows × 15 columns

df.shape

(32560, 15)

32560 rows x 15 columns

```
In [2]: df.shape
```

(32560, 15)

```
In [5]: df.columns # This will print the names of all columns
```

Index(['Age', 'Workclass', 'Fnlwgt', 'Education', 'Education\_num', 'Marital\_status', 'Occupation', 'Relationship', 'Race', 'Sex', 'Capital\_gain', 'Capital\_loss', 'Hours\_per\_week', 'Native\_country', 'Income'], dtype='object')

```
In [20]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32560 entries, 0 to 32559
Data columns (total 15 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   Age                 32560 non-null  int64
 1   Workclass           32560 non-null  object
 2   Fnlwgt              32560 non-null  int64
 3   Education            32560 non-null  object
 4   Education_num        32560 non-null  int64
 5   Marital_status       32560 non-null  object
 6   Occupation           32560 non-null  object
 7   Relationship         32560 non-null  object
 8   Race                 32560 non-null  object
 9   Sex                  32560 non-null  object
10   Capital_gain         32560 non-null  int64
11   Capital_loss         32560 non-null  int64
12   Hours_per_week       32560 non-null  int64
13   Native_country       32560 non-null  object
14   Income               32560 non-null  object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

```
In [7]: df.head() # Will give you first 5 recordsdf.head()
```

df.head() # Will give you first 5 recordsdf.head()																
Age	Workclass	Fnlwgt	Education	Education_num	Marital_status	Occupation	Relationship	Race	Sex	Capital_gain	Capital_loss	Hours_per_week	Native_country	Income		
0	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K	
1	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K	
2	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K	
3	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K	
4	37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	<=50K	

```
In [9]: df.tail(5)
```

	Age	Workclass	Fnlwgt	Education	Education_num	Marital_status	Occupation	Relationship	Race	Sex	Capital_gain	Capital_loss	Hours_per_week	Native_country	Income
32555	27	Private	257302	Assoc-acdm	12	Married-civ-spouse	Tech-support	Wife	White	Female	0	0	38	United-States	<=50K
32556	40	Private	154374	HS-grad	9	Married-civ-spouse	Machine-op-inspct	Husband	White	Male	0	0	40	United-States	>50K
32557	58	Private	151910	HS-grad	9	Widowed	Adm-clerical	Unmarried	White	Female	0	0	40	United-States	<=50K
32558	22	Private	201490	HS-grad	9	Never-married	Adm-clerical	Own-child	White	Male	0	0	20	United-States	<=50K
32559	52	Self-emp-inc	287927	HS-grad	9	Married-civ-spouse	Exec-managerial	Wife	White	Female	15024	0	40	United-States	>50K

## Data Preparation & Cleaning

```
In [17]: df.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32560 entries, 0 to 32559
Data columns (total 15 columns):
 # Column Non-Null Count Dtype
--- --
 0 Age 32560 non-null int64
 1 Workclass 32560 non-null object
 2 Fnlwgt 32560 non-null int64
 3 Education 32560 non-null object
 4 Education\_num 32560 non-null int64
 5 Marital\_status 32560 non-null object
 6 Occupation 32560 non-null object
 7 Relationship 32560 non-null object
 8 Race 32560 non-null object
 9 Sex 32560 non-null object
10 Capital\_gain 32560 non-null int64
11 Capital\_loss 32560 non-null int64
12 Hours\_per\_week 32560 non-null int64
13 Native\_country 32560 non-null object
14 Income 32560 non-null object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB

```
In [11]: df.describe()
```

	Age	Fnlwgt	Education_num	Capital_gain	Capital_loss	Hours_per_week
count	32560.000000	3.256000e+04	32560.000000	32560.000000	32560.000000	32560.000000
mean	38.581634	1.897818e+05	10.080590	1077.615172	87.306511	40.437469
std	13.640642	1.055498e+05	2.572709	7385.402999	402.966116	12.347618
min	17.000000	1.228500e+04	1.000000	0.000000	0.000000	1.000000
25%	28.000000	1.178315e+05	9.000000	0.000000	0.000000	40.000000
50%	37.000000	1.783630e+05	10.000000	0.000000	0.000000	40.000000
75%	48.000000	2.370545e+05	12.000000	0.000000	0.000000	45.000000
max	90.000000	1.484705e+06	16.000000	99999.000000	4356.000000	99.000000

## Checking the NULL values

```
In [13]: df.isnull()
```

	Age	Workclass	Fnlwgt	Education	Education_num	Marital_status	Occupation	Relationship	Race	Sex	Capital_gain	Capital_loss	Hours_per_week	Native_country	Income
0	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
32555	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
32556	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
32557	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
32558	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
32559	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False

32560 rows x 15 columns

```
In [15]: df.isnull().sum()
```

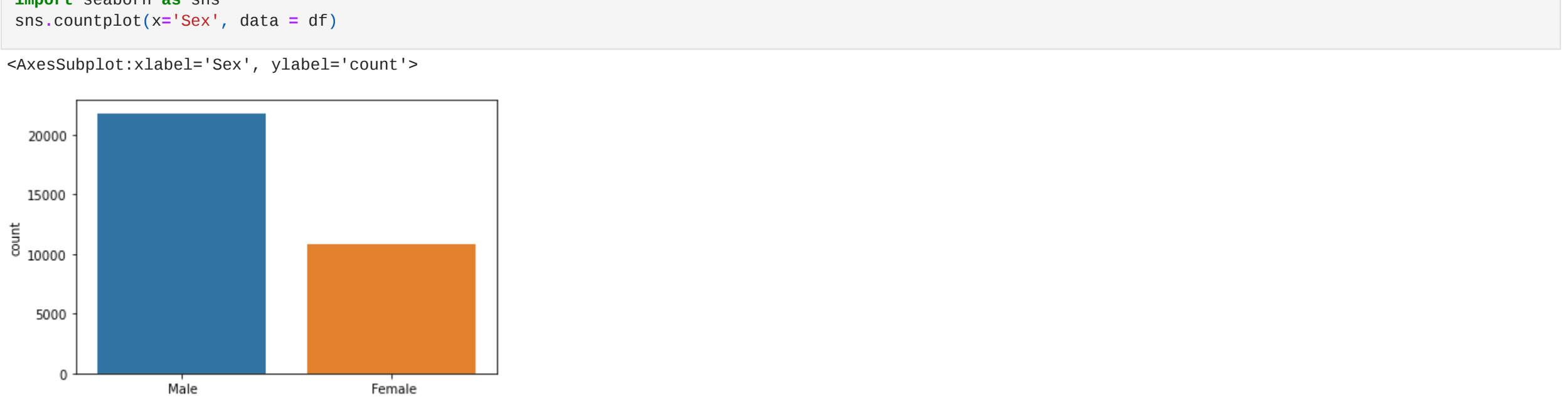
Age	0
Workclass	0
Fnlwgt	0
Education	0
Education_num	0
Marital_status	0
Occupation	0
Relationship	0
Race	0
Sex	0
Capital_gain	0
Capital_loss	0
Hours_per_week	0
Native_country	0
Income	0
dtype:	int64

```
In [16]: df.isnull().sum().sum()
```

0

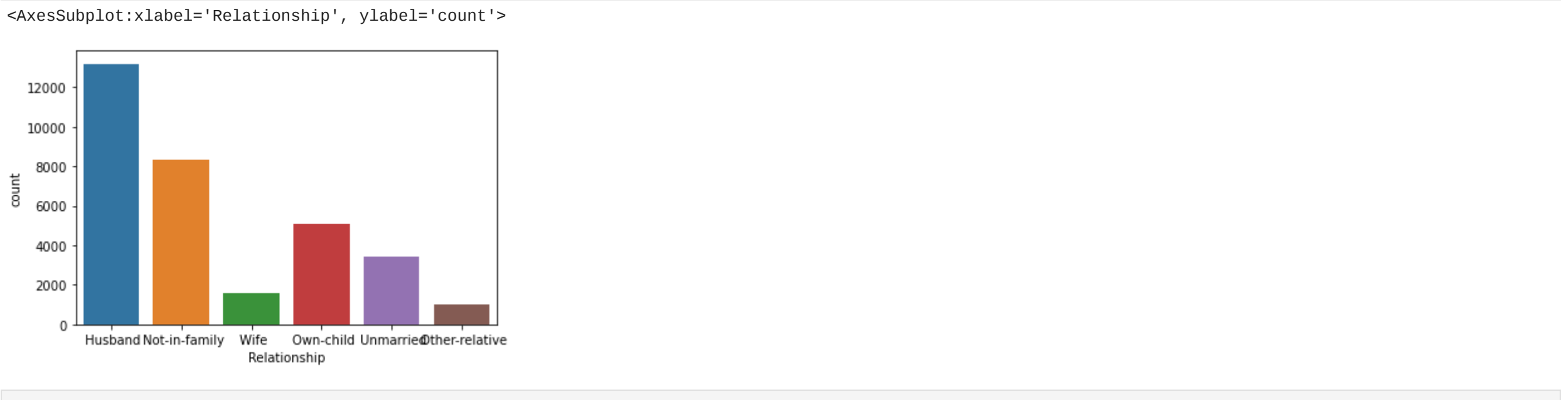
```
In [23]: import seaborn as sns
sns.countplot(x='Sex', data = df)
```

```
Out[23]: <AxesSubplot:xlabel='Sex', ylabel='count'>
```



```
In [24]: import seaborn as sns
sns.countplot(x='Relationship', data = df)
```

```
Out[24]: <AxesSubplot:xlabel='Relationship', ylabel='count'>
```



```
In [28]: import seaborn as sns
sns.countplot(x='Workclass', data = df)
```

```
Out[28]: <AxesSubplot:xlabel='Workclass', ylabel='count'>
```



```
In [32]: import seaborn as sns
sns.countplot(x='Income', data = df)
```

```
Out[32]: <AxesSubplot:xlabel='Income', ylabel='count'>
```



## Correlation

The correlation is a very important metric since it allows us to quantify the relationships between our variables, where we can see if they are directly or inversely related, and thus define which are the variables with most influence on the independent variables.

```
In [41]: corr_matrix=df.corr()
corr_matrix
```

	Age	Fnlwgt	Education_num	Capital_gain	Capital_loss	Hours_per_week
Age	1.000000	-0.076646	0.036527	0.077674	0.057775	0.068756
Fnlwgt	-0.076646	1.000000	-0.043159	0.000437	-0.010259	-0.018770
Education_num	0.036527	-0.043159	1.000000	0.122627	0.079932	0.148127
Capital_gain	0.077674	0.000437	0.122627	1.000000	-0.031614	0.078409
Capital_loss	0.057775	-0.010259	0.079932	-0.031614	1.000000	0.054256
Hours_per_week	0.068756	-0.018770	0.148127	0.078409	0.054256	1.000000

```
In [43]: sns.heatmap(df.corr(),annot=True)
```



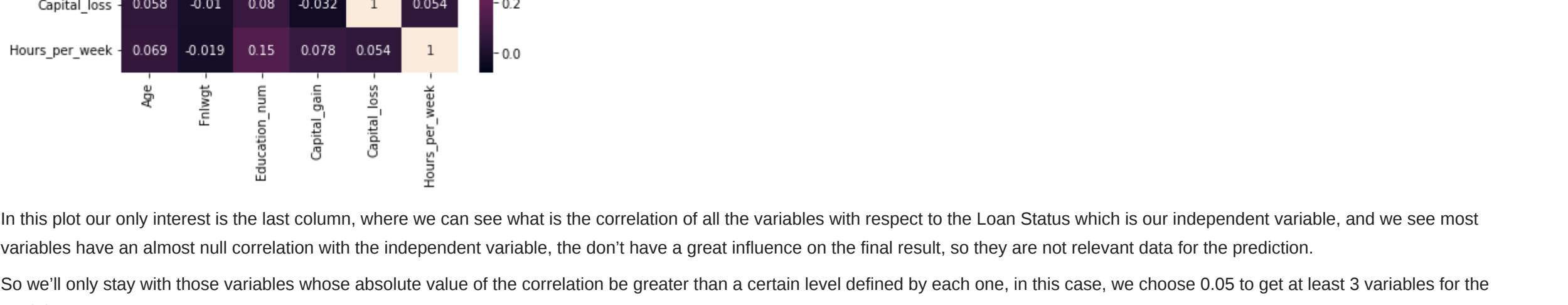
In this plot our only interest is the last column, where we can see what is the correlation of all the variables with respect to the Loan Status which is our independent variable, and we see most variables have an almost null correlation with the independent variable, the don't have a great influence on the final result, so they are not relevant data for the prediction.

So we'll only play with those variables whose absolute value of the correlation be greater than a certain level defined by each one, in this case, we choose 0.05 to get at least 3 variables for the model.

## Graphical Analysis

```
In [48]: import seaborn as sns
import matplotlib.pyplot as plt
sns.heatmap(df.isnull(),
plt.title('Null values'))
plt.show
```

```
Out[48]: <function matplotlib.pyplot.show(close=None, block=None)>
```



With the help of our work, we will be predicting the income of the population and analyzing the factors which strongly affect the income. We will be giving suggestions based on the result obtained which level of qualification can lead to a higher income and people of which age group are earning more. We will implement different models and find out which one is of higher accuracy. Also, we will consider the outputs of all the models and take a vote to determine the overall result.

## Conclusion

With the help of our work, we will be predicting the income of the population and analyzing the factors which strongly affect the income. We will be giving suggestions based on the result obtained which level of qualification can lead to a higher income and people of which age group are earning more. We will implement different models and find out which one is of higher accuracy. Also, we will consider the outputs of all the models and take a vote to determine the overall result.

## Thank you